

Comparative genomics of *Escherichia coli* causing  
bloodstream infections in the Hospital for Tropical  
Diseases, Ho Chi Minh city, Vietnam



**Tu Thi Phuong Le**

**Wellcome Trust Sanger Institute**

**University of Cambridge**

**Darwin College**

**This dissertation is submitted for the degree of Master of Philosophy in Biological Science**

**September 2017**

## **Declaration**

I hereby declare that:

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except where specifically indicated in the text and bibliography. My dissertation is not substantially the same as any that I have submitted for a degree or diploma or other qualification at any other University. I further state that no part of my dissertation has already been or is being concurrently submitted for any such degree, diploma or other qualification.

I confirmed that my dissertation does not exceed the limit of length of 20,000 words prescribed in the Special Regulations of the MPhil examination for which I am a candidate.

## Acknowledgements

This thesis was possible because of the collective efforts of many people, from the patient beds to lab benches and computer clusters, trying to deconvolute one of the oldest and cheekiest bacteria that cause big issues not only in a single referral hospital in Vietnam, but also across the globe.

This has been a tough but wonderful journey. I was extremely grateful to work on this project and have seen myself grow academically, intellectually, and personally in the past year in Cambridge.

First of all, I would like to thank my main supervisor Prof. Nicholas Robert Thomson for his unbelievable patience and guidance during the research. His continuous encouragement and support during difficult times is something that I always greatly appreciated.

I also would like to thank my former supervisor Prof. Stephen Baker from Oxford University Clinical Research Unit (OUCRU) for coming up with the idea for this project during the application process to the MPhil program as it turned out to be a very interesting project.

I would like to express my gratitude to my thesis committee members including Julian Rayner, Stephen Bentley, Estee Torok, and Kate Baker, who have provided valuable feedback and criticism.

I would like to thank members of Team 216 at the Sanger Institute for their tremendous daily support, as well as for providing me with a remarkably happy and friendly place to work. In no particular order, these include Angèle Bénard, Mathew Beale and Déborah Kreyenbihler. With no prior background on working with whole genome sequencing data, I therefore could not have completed the project without the guidance and patience of Alison Mather (Cambridge University), Daryl Domman, Eva Heinz and Alex Wailan (Sanger Institute), who never hesitating to spend time teaching me when I was stuck. I also would like to thank Ngure Kagia, another MSc student, who helped me to feel comfortable using Linux and computer clusters at the very beginning of my time at Sanger. Finally, I am indebted to Matthew Dorman for giving very constructive helpful feedback during my writing process.

The blood-derived *Escherichia coli* strains from this six-year period is the effort of all Microbiology staff at Hospital for Tropical Disease (Ho Chi Minh city, Vietnam), especially Nguyen Huu Hien and Tran Thi Ngan who carried out the blood culturing, identification, antimicrobial phenotype testing and sample archiving. The clinical metadata extracted from electronic hospital record database for this study was kindly provided by Nguyen Huu Hien, Dr Nguyen Phu Huong Lan and Dr Nguyen Van Vinh Chau. The second set of blood and stool isolates, as well as its associated clinical data, came from a big study on hospital acquired infections in the Intensive Care Unit lead by Dr Duong Bich Thuy, Mr James Campell, Dr Louise Thwaites, Prof Stephen Baker and Prof Guy Thwaites. I also would like to thank Nguyen Hoang Thu Trang, To Thi Thanh Huyen and Vong Vinh Phat for assisting me during the DNA extraction stage. I also would like to thank Dr Ho Dang Trung Nghia for discussions about the context of management and treatment of BSI patients in the hospital. My work relied on enormous sequencing

data of nearly 700 bacterial isolates, therefore I am indebted to the DNA core sequencing team and also the Pathogen Informatics team at the Sanger Institute for sample management, sequencing, and providing many analysis pipelines and bioinformatics supports. Many scripts and software I used for WGS analysis in this thesis were written by Simon Harris and Andrew Page.

Apart from working, I would like to give special thanks to Benjamin Bai, Simon Scutts, other friends and colleagues at Sanger and Darwin college badminton club for keeping me sane through many fun and joyful badminton matches.

Finally, this study could not have made possible without the support of the Wellcome Trust, who generously funded not only my university fees and living expenses, but also the research costs in both institutions in two nations, Sanger (UK) and OUCRU (Vietnam). I hope my study will shed a light on how the cheeky *E. coli* cause disease in Vietnamese patients as well as acquire drug resistance over time. Future work will be required to understand how to prevent and reduce the burden of disease in Vietnamese patients.

To my parents and little sister,  
for their endless and unconditional love.

## Table of Contents

Declaration.....	i
Acknowledgements.....	ii
Dedication.....	iv
Contents.....	v
Figures.....	vii
Tables.....	viii
Abstract.....	ix
<b>Chapter 1: Introduction.....</b>	<b>vii</b>
1.1. The history, taxonomy, and characteristics of <i>Escherichia coli</i> .....	1
1.1.1. Introduction .....	1
1.1.2. Commensal <i>E. coli</i> .....	1
1.1.3. The <i>E. coli</i> genome.....	2
1.1.4. Pathogenic <i>E. coli</i> .....	2
1.2. <i>E. coli</i> associated with bloodstream infections (BSIs), sepsis and bacteraemia .....	3
1.3. Hospital for Tropical Disease (HTD), Vietnam.....	4
1.4. Antimicrobial resistance in <i>E. coli</i> .....	4
1.5. AMR <i>E. coli</i> in Vietnam.....	5
1.5.1 ESBL – producing <i>E. coli</i> and carbapenemase- producing <i>E. coli</i> .....	6
1.5.2 Carbapenemase-producing <i>E. coli</i> .....	7
1.6. Mobile genetic elements (MGEs).....	7
1.7. Reservoirs of <i>E. coli</i> and transmission mechanism .....	8
1.8. Aims and objectives of the study .....	8
<b>Chapter 2: Comparative Genomics of BSIs caused by <i>E. coli</i> in the Hospital for Tropical Diseases in Vietnam .....</b>	<b>10</b>
2.1. Introduction .....	10
2.2. Methods.....	11
2.2.1. Strain collection .....	11
2.2.2 Ethics.....	11
2.2.3 Bacterial isolation .....	11
2.2.4 Metadata .....	12
2.2.5. Whole genome sequencing.....	12
2.2.6. Accessions .....	12
2.2.7. Genome assembly and annotation.....	12
2.2.8. Quality control of sequencing data.....	13
2.2.9. <i>In silico</i> MLST and phylogroup classification .....	13
2.2.10. Construction of core genome phylogeny .....	13
2.3. Results.....	14
2.3.1. Sequencing and QC .....	14
2.3.2. Basic patient demographic and clinical outcome of hospitalised patients with BSIs. ....	15
2.3.3. The genomic diversity and population structure of BSI and carriage <i>E. coli</i> from Vietnamese patients attending HTD.....	16
2.4. Discussion.....	25

<b>Chapter 3: The accessory genome of <i>Escherichia coli</i> of patients attending HTD. ....</b>	<b>27</b>
<b>3.1. Introduction .....</b>	<b>27</b>
<b>3.2. Methods.....</b>	<b>30</b>
3.2.1. <i>In silico</i> virulence, AMR, and plasmid replicon type gene detection.....	30
3.2.2. Antimicrobial phenotype testing .....	30
3.2.3. Statistical analysis.....	30
<b>3.3. Results .....</b>	<b>31</b>
3.3.1. The distribution of virulence factors between blood and rectal swab isolates .....	31
3.3.2. The distribution of virulence factors between different BAPS lineages/ phylogroups/STs.....	<b>33</b>
3.3.3. Antimicrobial resistance phenotype between blood/rectal swab isolates .....	41
3.3.4. Antimicrobial resistance genes distribution in blood/ rectal swab .....	42
<b>3.4. Discussion.....</b>	<b>50</b>
<b>Chapter 4: Understanding the relationships between BSI and carriage <i>E. coli</i> isolates in patients attending HTD, Vietnam. ....</b>	<b>53</b>
<b>4.1. Introduction .....</b>	<b>53</b>
<b>4.2. Methods.....</b>	<b>54</b>
4.2.1 Study participants.....	54
4.2.2. Isolate collection.....	54
4.2.3. Phylogenetic analysis .....	54
4.2.4. AMR genes, replicon types and virulence genes identification.....	54
<b>4.3. Results .....</b>	<b>54</b>
4.3.1. Patient and sample recruitment .....	54
4.3.2. Genomic diversity among rectal swab and blood derived <i>E. coli</i> isolates on admission to HTD ICU .....	55
4.3.2. Longitudinal diversity .....	59
<b>4.4. Discussion.....</b>	<b>63</b>
<b>Chapter 5 Final discussion.....</b>	<b>65</b>

## Figures

<b>Figure 2.1.</b> Sequencing and QC work flow for the samples included for whole genome sequencing (WGS) in this study .....	14
<b>Figure 2.2:</b> Comorbidity and outcome in patients with <i>E. coli</i> BSIs.....	16
<b>Figure 2.3:</b> The pan genome of 687 <i>Escherichia coli</i> .....	18
<b>Figure 2.4:</b> Population structure of 665 <i>E. coli</i> isolates analysed in this study.....	19
<b>Figure 2.5:</b> The phylogenies of <i>E. coli</i> taken from HTD, including reference strains, reported alongside clinical metadata for <i>E. coli</i> causing BSI.....	22
<b>Figure 2.6:</b> STs distribution of blood and rectal swab isolates.....	23
<b>Figure 2.7:</b> A core genome tree visualised in Phandango for blood and rectal swab isolates.....	24
<b>Figure 3.1:</b> Volcano plot for odds ratio versus $-\log_{10}$ p-value of gene distribution between invasive and carriage isolates.....	33
<b>Figure 3.2:</b> Phylogeny of <i>E. coli</i> lineages and virulence genes found in each isolate. ....	34
<b>Figure 3.3:</b> The distribution of AMR resistance between blood and carriage isolate.....	42
<b>Figure 3.4:</b> Total number of resistance genes and replicon types.....	43
<b>Figure 3.5:</b> AMR genotypes and phenotypes across all <i>E. coli</i> isolates.....	44
<b>Figure 3.6:</b> The genetic environment of <i>mcr-1</i> on a chromosomal contig.....	45
<b>Figure 4.1:</b> Phylogenetic diversity among rectal swab and blood derived <i>E. coli</i> isolates from HTD .....	57
<b>Figure 4.2:</b> Pairwise SNP distances within and between patients based on the variant calls from the core gene alignment.....	58
<b>Figure 4.3:</b> Twenty longitudinal isolates and from seven patients .....	59
<b>Figure 4.4:</b> Graph shows phylogenetic relationship between isolates in patient P3 (03-283).....	63



## Tables

<b>Table 2.1:</b> Clinical characteristic of <i>E. coli</i> BSIs patients.....	<b>15</b>
<b>Table 3.1:</b> Common virulence factors identified in ExPEC.....	<b>28</b>
<b>Table 3.2:</b> The percentage of isolates carrying specific virulence genes in the 12 most dominant sequence types of <i>E. coli</i> found in this study.....	<b>36</b>
<b>Table 3.3:</b> <i>bla</i> CTX-M distribution among different STs .....	<b>46</b>
<b>Table 3.4:</b> Plasmid replicon type identified in <i>E. coli</i> collection.....	<b>47</b>

## Abstract

*Escherichia coli* (*E. coli*) is a versatile bacterium, with the capability to act not only as a commensal coloniser but also as a pathogen that can cause invasive disease. Currently, *E. coli* is the leading cause of bloodstream infections in both developed and developing countries, accounting for 25-30% of bacteraemia cases globally. *E. coli* bacteraemia is further exacerbated by the emergence of antimicrobial resistance, particularly extended spectrum  $\beta$ -lactamases (ESBLs) and carbapenemases in Gram-negative bacteria. Understanding the nature and diversity of *E. coli* causing bloodstream infections is crucial for the enhancement of infection control measures, to minimise the further emergence of antimicrobial resistant isolates, and for the reduction of morbidity and mortality.

We used whole-genome sequencing to analyse the population structure of 506 invasive and 159 carriage *E. coli* isolates collected at The Hospital for Tropical Diseases, Vietnam during 2010-2015, and to look for genetic signatures that differentiate them from one another. We found substantial diversity therein, both in the total number of sequence types present, and in the number of resistant genes carried. Among blood isolates, ST131, ST95, ST69, ST1193 and ST73 were dominant STs, while from the rectal swab (carriage) isolates, ST131, ST1193 and ST648 were most prevalent. All of these STs from blood and rectal swab samples remained dominant over the entire study period, the exception being ST1193. Interestingly ST1193 was not present in Vietnam prior to 2011, but once introduced, it quickly emerged and replaced other more drug sensitive *E. coli* clones. From longitudinal sampling and paired *E. coli* isolated from blood and rectal swab from the same patients, we demonstrated that majority of *E. coli* infections (57%) were acquired from the patients' own gut microbiota. This study shows that several genetic factors, including genes that mediate adhesion (fimbriae and pili), iron acquisition (siderophores), immune evasion (capsule synthesis) and toxin elaboration (haemolysin and cytotoxic necrotizing factor 1) are significantly more common in invasive *E. coli* than in carriage strains. We were also able to begin to explain anomalies in the patterns of antimicrobial resistance of blood and rectal swab *E. coli* isolates using the population structure we defined here.

Taken together, we have generated a genetic framework for future studies focusing on *E. coli* BSI at HTD. These data also demonstrate that a combination of both virulence genes and antimicrobial resistance genes are essential for the success of certain lineages in causing invasive disease in Vietnamese patients. Perhaps more importantly we also show that patients who develop bacteraemia in hospital did so largely through infection with isolates already present in their intestinal tract and so were predominantly community acquired invasive infections rather than hospital acquired infections, contrary to our original hypothesis.